

## The Visual Storytelling Task

- Given a **photo stream**, the machine is trained to **generate a coherent story** in natural language to describe the photos.

### Example: visual storytelling



it was a great time at the wedding. There was a lot of the wedding. It was a great time at the reception. It was a great time to the end of the day. At the end of the boat, we went to the boat.

- Issue:** previous methods focused on optimizing automatic metrics which are not especially designed for text quality

## ReCo-RL: Composite Rewards

- Our approach:** design a reward function that correlates with story quality
- Specifically, use three different metrics:

**Relevance**  
Extract important caption entities that appears in the image

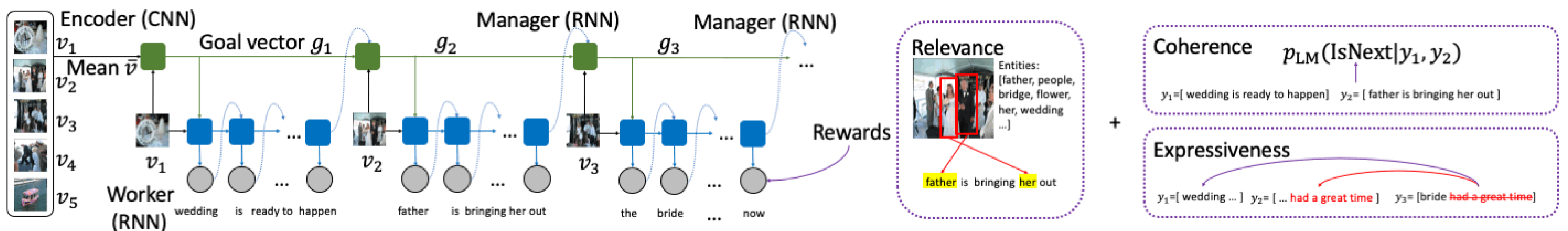
**Coherence**  
Estimate whether two sentences are coherent by a next sentence predictor

**Expressiveness**  
Penalize repeated patterns by comparing w/ previous sentences

- Train the model by REINFORCE and MLE

$$\theta \leftarrow \theta + \eta_1 \frac{\partial J_{MLE}(\theta, \mathcal{D}')}{\partial \theta} + \eta_2 \frac{\partial J_{RL}(\theta, \mathcal{D}')}{\partial \theta}$$

## Model Architecture: Manager-Worker w/ Three Rewards



## How to evaluate the text quality? Human Evaluation & Automatic Metrics

- Baselines:** MLE; AREL; HSRL; ReCo-RL
- Automatic Metric:** METEOR; BLEU; CIDEr; SPICE; ROUGH-L
  - Our MLE/BLEU-RL/ReCo-RL are competitive to SOTA
  - Automatic metrics are not designed for visual storytelling
- Human evaluation:** relevance, coherence, expressiveness in Table 2.
  - ReCo-RL outperforms other baselines
  - Student's paired t-test with  $\rho < 0.05$
  - 862 turkers show substantial agreement
- ReCo-RL gets higher reward scores on test set in Table 3.

Method	METEOR	ROUGE	CIDEr	BLEU-4	SPICE
AREL	35.2	29.3	9.1	13.6	8.9
HSRL	30.1	25.1	5.9	9.8	7.5
MLE	34.8	30.0	7.2	14.3	8.5
BLEU-RL	35.2	30.1	6.7	14.4	8.3
ReCo-RL	33.9	29.9	8.6	12.4	8.3

Table 1: Comparison between different models on METEOR, ROUGE-L, CIDEr, BLEU-4 and SPICE.

Method	Relevance	Coherence	Expressiveness
HSRL	1.95	7.21	33.27
AREL	3.27	9.90	34.98
MLE	5.46	7.92	30.76
BLEU-RL	2.17	12.40	30.41
ReCo-RL	10.39	12.74	39.37

Table 3: Comparison between different models on three rewards, i.e., Relevance, Coherence and Expressiveness.

Aspects	AREL	ReCo-RL	Tie	Agree	HSRL	ReCo-RL	Tie	Agree	MLE	ReCo-RL	Tie	Agree	BLEU-RL	ReCo-RL	Tie	Agree
R	27.6%	62.2%	10.2%	0.72	36.1%	53.8%	10.1%	0.74	27.0%	64.1%	8.9%	0.49	17.6%	74.5%	7.9%	0.78
C	31.3%	58.7%	10.0%	0.78	38.0%	51.9%	10.1%	0.80	34.3%	57.7%	8.0%	0.53	18.9%	72.3%	8.8%	0.71
E	32.4%	58.6%	9.0%	0.68	38.6%	53.3%	8.1%	0.72	30.5%	61.0%	8.5%	0.55	19.5%	71.5%	9.0%	0.62

Table 2: Pairwise human comparison between ReCo-RL and three methods on Relevance, Coherence and Expressiveness.

## Qualitative Analysis: what does the generated story look like?

- BLEU-RL and HSRL generate uninformative segments.
- AREL generates topically incoherent stories.
- ReCo-RL generates rare entities such as "sign" and "flags"
- Stories generated by ReCo-RL obtains higher reward scores.

Method	Quality Metrics	Quality Metrics			
		R	C	E	B
BLEU-RL	a group of friends gathered together for dinner. the turkey was delicious. the guests were having a great time. at the end of the night, we had a great time. at the end of the night, we had a great time.	2.47	11.06	37.10	73.57
ReCo-RL	a group of friends gathered together for a party. the turkey was delicious. it was a delicious meal. everyone was having a great time. after the party, we all sat down and talked about the night. my friend and [female] were very happy to drink.	3.32	16.99	41.71	78.51

Methods	Generated Story	Generated Story
AREL	the officers of the military officers are in charge of the military. he was very proud of his speech. the meeting was a great success. the president of the company gave a speech to the audience. we had a great time.	the wedding was held at a church. the wedding was beautiful. the bride and groom cut the cake. the bride and groom were very happy. the whole family was there to celebrate.
HSRL	i was so excited to see my new team. he was very happy to see the new professor. i had a lot of time to talk about. i had a great time. i had a great time.	the wedding was a beautiful wedding. the bride and groom cut the dance together. the bride and groom were very happy to be married. the bride and groom were very happy. at the end of the night, the bride and groom were happy to be married.
BLEU-RL	at the end of the day, the men were very proud of the military. they had a lot of people there. this is a picture of the meeting. a group of people had a great time. after the end of the day, we all had a lot of questions.	it was a beautiful day for the wedding. at the end of the night, the bride and groom were very happy. the bride and groom were very happy. she was so happy to be married. the bride and groom pose for pictures.
ReCo-RL	today was a picture of the military officer. he was ready to go to the organization. they were very happy to see the awards ceremony. the speaker was very excited to be able to talk about the meeting. everyone was having a great time to get together for the event after the ceremony. we all had a lot of people there.	it was a beautiful day at the wedding party. the bride and groom were so happy to be married. [female] was happy and she was so excited to celebrate. she had a great time to take a picture of her wedding. all of the girls posed for pictures.